

BEHAVIOR RECOGNITION SYSTEM

Reference to Prior Applications

This application claims priority of U.S. provisional application Serial No. 60/127,510 filed April 2, 1999, and is a continuation-in-part of U.S. patent application Serial No. 09/371,460, filed August 10, 1999, the entire contents of both of which are
5 incorporated herein by reference.

Field of the Invention

This invention relates generally to automated image recognition and, in particular, to computer-based behavior recognition methods and apparatus.

Background of the Invention

10 Differentiating between normal human activity and suspicious behavior is a difficult task, whether performed by a sensing device or a human observer. A human observer would find such a task tedious and costly to perform in money and time. Fortunately, a sensing system is not bothered by 'tedious' tasks. Such a system could be implemented to prune out obviously normal behavior, and tag human activities which
15 could be suspicious and would therefore need more attention by a human operator.

However, such "behavior recognition" systems have not been developed due to the difficulty of identifying and classifying such motions. Consider an urban environment

as illustrated in Figure 1. In such an environment, there are many moving objects and people, most of which are not actively engaged in criminal or destructive behavior. All of these activities would have to be viewed and identified before they could be removed from potential further consideration by a human.

5

Summary of the Invention

The present invention is directed to the automated classification and identification of human activities. The requisite heuristics involve recognition of information-bearing features in the environment, and the determination of how those features relate to each other over time. The approach resides in a behavior recognition sensor system which
10 identifies simultaneously created gestures using a gesture recognition system. The underlying gesture recognition system performs the task of determining the behavior (state) of objects in motion.

According to the invention, general kinematic relationships for various bodies (human, animal, robotic) are modeled. Next, specific link combinations for each are
15 parameterized and modeled. This enables the system to recognize such motions as the various human motion gaits (if the links are legs), the throwing of an object (if the links are arms), or any other type of motion and orientation. A whole-body (human, vehicle, or other) link representation and dynamic model is then developed.

The state of an object (a human or a device) can take three forms.

- The object is static (motionless) and cannot be easily moved by a human agent. Such objects would include walls, filing cabinets, telephone poles, and the like. These are identified and localized with respect to the vision system;
- The object is static but can be manipulated by a human agent. Such objects would include weapons, chairs, etc. These are localized with respect to the vision system, tagged, and identified and evaluated with regard to their potential use by human agents; and
- The object is in motion. Such objects include self-mobile objects such as humans and vehicles, as well as objects that are carried by humans or vehicles, such as guns or briefcases.

Whenever an object is identified, the system preferably first determines which state the object is in based upon gesture recognition. The term "gesture" not only refers to the (dynamic) motion of an object, but also to the state of an object that is not moving (static). For example, a gun being pointed by a human should definitely be classified as a gesture. Therefore, determining behaviors also means identifying combinations of static and dynamic gestures.

An object is composed of many connected non-deformable links. For example, a person is made up of a torso (one link), a head (one link), two arms (each with two links), two hands (one link each, not counting the fingers), two legs (each with two links), and two feet (one link each). Each link of an object has its own gesture motion, which, when examined in relation to the other links in the object, can be used to determine the overall

state of the system. For example, when just examining the gestures created by the foot and knee joint of a human, one can see that the motions of those features are different depending on whether a person is walking or running. Even when not in motion, the relationship of those non-deformable links gives rise to information about the object's
5 state. A person who is standing still will have a different kinematic link relationship to one who is sitting, or lying down.

Such gross human torso features can be readily identified (using a model that matches general body features) with the vision system described in co-pending U.S. Patent Application Serial No. 09/371,460. Building on this technology, the gesture
10 recognition module of this invention determines if a dynamic motion is occurring, and uses that information with kinematic link relationships, develop a hypothesis about the overall state of the object or objects in the field of view.

In a sense, then, the invention provides an automatic method for classifying and categorizing such dynamic motions and gaits. Such state recognition is not limited to
15 humans (and other animals), however. A vehicle's state can also be determined by examining the various moving parts, such as the body and the tire motions. Even unknown devices (such as mobile robots) can be classified by examining their motion and behavioral characteristics.

Although subject visualization is preferred, alternative approaches to
20 gait/behavior recognition may be employed according to the invention, including electro-mechanical methods of tracking the body to identify human motion. This includes using:

- electrogoniometers and electrogoniometric systems (EGM)
- passive reflective and actively illuminated markers to calculate positions and velocities, i.e. raw data.
- force data gathered from a force plate or force dynamometer. Used to calculate
5 internal joint moments causing motion.
- strain gage or piezoelectric transducers to measure ground reaction forces.

Brief Description of the Drawings

FIGURE 1 is an overview of potential behavior recognition system uses;

10 FIGURE 2 is a drawing of a gesture recognition system according to the invention;

FIGURE 3 is a gesture recognition system flow chart;

FIGURE 4 is a signal flow diagram of a gesture recognition system according to the invention;

15 FIGURE 5 is a drawing which shows example gestures in two dimensions;

FIGURE 6 shows three example gestures;

FIGURE 7 is an example of a 24-gesture lexicon according to the invention;

FIGURE 8 depicts a Slow-Down gesture;

FIGURE 9 depicts a Move gesture;

FIGURE 10 depicts an Attention gesture;

FIGURE 11 depicts a Stop gesture;

FIGURE 12 shows Right/Left Turn gestures;

5 FIGURE 13 shows an "Okay" gesture;

FIGURE 14 shows a Freeze gesture;

FIGURE 15 provides three plots of a human created one dimensional X-Line
oscillating motion;

FIGURE 16 shows possible lines associated with $x(t,p)=p_0+p_1t$ and their
10 equivalent representation in the p-parameter space;

FIGURE 17 illustrates parameter fitting wherein a rule is used for q to bring the
error to zero;

FIGURE 18 plots different (x_i, y_i) data points resulting in a different best fitting q
line;

15 FIGURE 19 depicts a recursive linear least squares method for updating q with
subsequent (x_i, y_i) data points;

FIGURE 20 illustrates an algorithm for determining a specific gesture model according to the invention;

FIGURE 21 is an exaggerated representation of a residual error measurement;

FIGURE 22 is a plot which shows worst case residual ratios for each gesture
5 model, wherein the lower the ratio, the better the model;

FIGURE 23 illustrates how two perpendicular oscillatory line motions may be combined into a circular gesture;

FIGURE 24 shows how a bounding box may be placed around a hand associated with a gesture;

10 FIGURE 25 provides descriptions from the bounding box of Figure 24;

FIGURE 26 shows example gestures;

FIGURE 27 is a schematic of hand-tracking system hardware according to the invention;

FIGURE 28 is a flowchart of a color tracking system (CTS) according to the
15 invention;

FIGURE 29 depicts a preferred graphical user interface of the CTS;

FIGURE 30 illustrates the application of target center from difference image techniques;

FIGURE 31 illustrates a color matching technique;

FIGURE 32 is a representation of an identification module;

FIGURE 33 is a simplified diagram of a dynamic gesture prediction module according to the invention; and

- 5 FIGURE 34 is a simplified diagram of a behavior recognition module according to the invention.

Detailed Description of the Invention

- The core of the behavior recognition system is a set of dynamic and static gesture recognition modules. This section details the overall gesture recognition system used by
10 the behavior recognition system.

The Gesture Recognition System

- Figure 2 presents a system overview of a gesture controlled self service machine system. Figure 3 shows a flow chart representation of how a vision system views the
15 gesture created, with the image data sent to the gesture recognition module, translated into a response, and then used to control a SSM, including the display of data, a virtual environment, and devices. The gesture recognition system takes the feature positions of the moving body parts (two or three dimensional space coordinates, plus a time stamp) as

the input as quickly as vision system can output the data and outputs what gesture (if any) was recognized, again at the same rate as the vision system outputs data.

The specific components of the gesture recognition system is detailed in Figure 4, and is composed of five modules:

- 5 • G: Gesture Generation
- S: Sensing (vision)
- I: Identification Module
- T: Transformation
- R: Response

10 At a high level, the flow of the system is as follows. Within the field of view of one or more standard video cameras, a gesture is made by a person or device. During the gesture making process, a video capture card is capturing images, producing image data along with timing information. As the image data is produced, they are run through a feature tracking algorithm which outputs position and time information. This position
15 information is processed by static and dynamic gesture recognition algorithms. When the gesture is recognized, a command message corresponding to that gesture type is sent to the device to be controlled, which then performs an appropriate response. The five modules are detailed below.

Gesture Creator

20 In the Gesture Creator module, a human or devices creates a spatial motion to be recognized by the sensor module. If one camera is used, then the motion generated

is two dimensional and parallel to the image plane of the monocular vision system. For three dimensional tracking (as is also done with this system), stereo vision using two or more cameras are used.

This gesture recognition system is designed to recognize consistent yet non-perfect motion gestures and non-moving static gestures. Therefore, a human can create such gestures, as well as an actuated mechanism which could repeatedly create perfect gestures. Human gestures are more difficult to recognize due to the wide range of motions that humans recognize as the same gesture. We designed our gesture recognition system to recognize simple Lissagous gesture motions (repeating circles and lines), repeated complex motions (such as “come here” and “go away quickly” back and forth hand motions which we define as “skew” gestures), and static hand symbols (such as “thumbs-up”).

With regards to human generated gestures used for communication or device control, we chose gestures to be identified based on the following:

- Humans should be able to make the gestures easily.
- The gestures should be easily represented mathematically.
- The lexicon should match useful gestures found in real world environments.

For the dynamic (circular and skew) gestures, these consist of one-dimensional oscillations, performed simultaneously in two or three dimensions. A circle is such a motion, created by combining repeating motions in two dimensions that have the same magnitude and frequency of oscillation, but with the individual motions ninety degrees

out of phase. A “diagonal” line is another such motion. We have defined three distinct circular gestures in terms of their frequency rates: slow, medium, and fast. An example set of such gestures is shown in Figure 5. These gestures can also be performed in three dimensions, and such more complex motions can be identified by this system.

- 5 The dynamic gestures are represented by a second order equation, one for each axis:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2\end{aligned}$$

- More complex second order models are used to recognize more complex gestures
10 (discussed later). This gesture model has no “size” parameter. θ_1 is a frequency measure, and θ_2 is a drift component. The gestures were named “large”, “small”, “fast”, and “slow” due to the human motions used to determine the parameters (see Figure 6). A fast small circle is used to represent a fast oscillation because humans cannot make fast oscillations using large circles.

- 15 For example, a total of 24 gestures are possible when the following are distinct gestures: clockwise and counter-clockwise circles, diagonal lines, one dimensional lines, and small and large circles and lines. Geometric constraints are required to expand the lexicon, because different gestures can result in the same parameters. Figure 5 shows motions that would cause an identifier to produce the same frequency measure and drift
20 components as it would produce when identifying a slow large circle. When x and y

oscillating motions are 90 degrees out of phase, a clockwise circle is produced. Motions that are 270 degrees out of phase result in a counter clockwise circle. In phase motions produce a line with a positive slope. When the motions are 180 degrees out of phase, a line with a negative slope is produced. We can create additional gestures from the fast
5 small circle in the same manner.

As with the previous gestures, additional gestures can be created from these two gestures by varying the phase relationships. Figure 7 shows a representation of the 24 gestures in possible lexicon. Even more gestures are possible when the third dimension is used.

10 Phase relationships are determined as follows. During the gesture, the x's and y's (and z's, if the system is set up for three dimensions) minimum and maximum image plane positions are computed. If the x and y motions are out of phase, as in a circle, then when x or y is minimum or maximum, the other axis's velocity is large. The direction (clockwiseness in two dimensions) of the motion is determined by looking at the sign of
15 this velocity component. Similarly, if the x and y motion are in phase, then at these extremum points both velocities are small.

Example dynamic gestures used for real world situations are taken from a 1987 Army Training Manual. A "Slow Down" gesture is a small x-line created to one side of the body (Figure 8, left side). A "Day Move" gesture is a counterclockwise large slow
20 circle (Figure 9, left side). The "Attention" gesture is a large y-line overhead motion

(Figure 10). These three gestures are representative of the motion gestures used throughout the Army manual.

Static gestures are represented as geometric templates. Four gestures are shown and are representative of the static gestures which can be represented and identified by this gesture recognition system. Additionally, language gestures, such as American Sign Language gestures, can also be recognized.

The example static gestures are:

- Halt - stop hand above head (Figure 11 - left side of figure).
- Left and Right turn - fingers together, palm out, facing left or right (Figure 12 - left side of figure).
- Message Acknowledge (OK) - thumb up (Figure 13).
- Freeze - Fist at head level (Figure 14).

Identifying Moving Gestures Represented as a Dynamic System

The gesture recognition system identifies a moving gesture by its dynamics - that is, the structure of its positions in space over time. The system translates the motion information into parameters which are used to develop commands for controlling data outputs and actuated mechanisms. For example, the speed at which a person waves a robot away might directly affect a robot arm's velocity or a mobile robot's speed. In order for recognition to occur, a representation for human gestures is required, from which a computational method for determining and recognizing specific gestures can be derived.

Although we make these gestures in two and three dimensions, the explanation now detailed is described simply dimension as a basic one dimensional gesture as a simple example to clarify the distinction between the “shape” and the “dynamics” of a gesture. The techniques for identifying this basic gesture will be used to identify similar
5 oscillatory motions occurring in two and three dimensions.

First, a dynamic system gesture representation is determined, both the model for representing the oscillatory gestures and parameter determination scheme was developed. For this system a Linear Least Squares method was an on-line computationally efficient technique which allowed us to use a linear-in-parameters gesture model.

10 The representative planar gesture used throughout this section to exemplify our method consists of a family of oscillating motions which form a (roughly) horizontal line segment (“x-line motion”). As discussed earlier, a human is incapable of reliably generating a perfect sinusoidal motion. Figure 15 illustrates the imperfections of a human created x-line motion viewed in three plots. The plots represent the position of the
15 gesture over time, $x(t)$. Viewing position with respect to time in contrast to position and velocity over time provides insight into how we propose to represent gestures. Plot A shows the planar motion in x-position and y-position coordinates, with the gesture's motion constrained to the x-axis. Thus, the “shape” of the motion conveys relatively little information. Plot B shows the same gesture in x-position plotted against time,
20 emphasizing the oscillatory behavior we wish to capture. Plot C represents the record of x-velocity plotted against x-position over time. We will find it most convenient to

represent this motion as it evolves over time in this position versus velocity space, which is called the “phase plane”. Of course, when a human creates a gesture, the resulting motion does not translate into the perfect sinusoid of plot B or a perfect circle of plot C. Instead, there is a natural range of variation that we would nevertheless like to associate
5 with the same gesture. This association we find most naturally achievable in phase space.

For this dynamic gesture recognition module, a computationally effective mathematical representation for the gesture plotted in Figure 15 is required. A general representation for time functions might take the form

$$x(t) = ?,$$

10 where “?” would be replaced with some structure based on measurable features which are used to classify the gesture. Of course, there are an infinite number of possible measurable features.

We can make the number of classifications (the “feature space” dimension) finite by restricting the form of the representations. Instead of representing gestures as $x(t)$, the
15 representation might be constrained through the use of a *parameter* vector, resulting in $x(t,p)$. The feature space dimension is then equivalent to the number of parameters we store. For example, when:

$$x(t,p)=p_0+p_1t, t>0,$$

the only possible gestures that we can represent are lines described by the two parameters
20 slope, p_1 , and intercept p_0 (see Figure 16).

Even with a finite dimensional representation, each unique motion is represented by its own distinct parameters. However, our intuition about human gestures tells us that certain distinct motions should have the same classification. Consider the x-line oscillating gesture discussed earlier. Whether the gesture starts at the left side of the line or the right side (for example, $x(0)=-1$ or $x(0)=+1$), the resulting motions would still be identified by a human as the same gesture. Therefore, another type of representation seems desirable.

Since a human hand forms a gesture, we could imagine a representation in terms of the force exerted by the person's arm muscles. Alternatively, we might imagine representing the gesture as a function of the nerve impulses that travel from the brain to the arm's muscles. However, quite clearly, most of the countless types of such "internal" representations are presently impossible to quantify in any useful manner. Four hundred years ago, Newton developed a parsimonious representation of physical motions based on their dynamic properties,

$$\dot{x}(t) = f(x)$$

A dynamic system is a mathematical model describing the evolution of all possible states in some state space as a function of time [Hirsch 74] [Arnold 78]. The set of all possible states is a state space. Given an initial state, the set of all subsequent states as it evolves over time is a "trajectory" or "motion". For any initial condition, the future evolution of the states in a trajectory remains within that trajectory (the trajectory is an invariant set). Thus, all that is required to describe a particular spatial motion is the differential equation

representation and its initial conditions. We use a deterministic representation, as opposed to a stochastic one, because we believe these oscillatory motions are best represented by sine waves or a sum of exponentials as opposed to characteristics based on statistical properties. As with the geometric representation, there are an infinite number of gesture classifications of the form $\dot{x}(t) = f(x)$. However, as before, we can choose a vector of tunable parameters to make the number of gesture classifications finite. Such representation has the form:

$$\dot{x}(t) = f(x, \theta)$$

where θ represents the tunable parameters. Fixing the value of θ in a given representation yields a unique set of motions, with different initial conditions, described by $\dot{x}(t) = f(x, \theta)$. Motivated by the way humans interpret gestures, we associate an entire set of motions with one specific gesture. Thus, choosing different values of θ in a given representation results in a “family” of trajectories sets - a “gesture family.” For example, consider a oscillatory line gesture, the motion of which is constrained to the x-axis. This gesture can be represented in the following two dimensional state space:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1\end{aligned}$$

where x_1 represents the position of the gesture, x_2 is its velocity, and θ_1 is a specified negative parameter. For any constant $\theta > 0$, all trajectories satisfy $-\theta_1 x_1^2 + x_2^2 = \text{const}$ as can be seen by direct differentiation.

We conceive of a specific gesture as a family of sets of trajectories. Referring to the gesture in Figure 15, a human can start the gesture at any point (initial condition) in its trajectory, and the gesture should still be identified as the same oscillating line.

We will now represent a given family of gestures (family of sets of trajectories) by
5 a mathematical model which contains a finite number of tunable parameters. A mathematical model described by differential equations, as above, allows the development of a computational scheme that will determine which parameters, the values of θ_i 's, correspond to a specific gesture. The set of all valid parameters is the parameter space. The parameter space defines the family of gestures which can be represented by
10 the model. In order to categorize a finite number of gestures in this family and to permit further variability in the exact motions associated with a particular gesture within this family, we partition the parameter space into a finite number of cells - the "lexicon" - and associate all the parameter values in the same cell with one gesture.

We will "invent" certain differential equations, composed of state variables and
15 parameters, which intuition suggests may represent human gestures. Such differential equation models can be divided into two types: non-linear-in-parameters (NLIP) and linear-in-parameters (LIP). The two models can be further subdivided into linear-in-state (LIS) and non-linear-in-state (NLIS). It is advantageous to use a NLIP (with NLIS) model because it covers, by definition, a much broader range of systems than an LIP model.
20 However, for reasons to be discussed below, we find it expedient to use a LIP model for our gesture representation. We have chosen to represent planar oscillatory gestures as a

second order system with the intuition that a model based on the acceleration behavior (physical dynamics) of a system is sufficient to characterize the oscillatory gestures in which we are interested. This system's states are position and velocity. However, the vision system we use to sense gestures yields only position information. Since velocity is not directly measured, then either the parameter identification method could be combined with a technique for observing the velocity, or the velocity could be determined through position differences. In this section we show techniques for determining gesture parameters both when the velocity state is observed, and when it is obtained through position differences. By examining the utility of each technique, we develop an appropriate form of the gesture model and parameter identification method.

A difficulty with using human created gestures is that the underlying true physical model is unknown. Also, because people cannot precisely recreate even a simple circular gesture, multiple sets of parameters could represent the same gesture. Simulations are used both to determine a viable gesture model and to determine if it is possible to discover appropriate parameters for each gesture despite variations in motion.

We had chosen to represent motion gestures using dynamic systems. Our next task was to determine a model and a method for computing the model's parameters such that the model's parameters will best match an observed gesture motion. Figure 17 illustrates how the gesture's position is used as an input, with $\hat{\theta}$ representing the unknown parameter values that we wish to match with the "true" parameter values, θ . If these

values match, then the error between the true states \dot{x} and the observed states $\hat{\dot{x}}$ will go to zero.

Our choice of a model and parameter determination scheme was based on an exploration of the following issues:

- 5 • Off-line batch techniques versus on-line sequential techniques. We desire our gesture recognition system to identify gestures *as they are generated*, which requires an on-line technique. Also, the measure of how well a motion matches a gesture's parameters needs to be updated “on-line”.
- 10 • State availability. Using a vision system to sense gestures results in image plane position information. However, we are using a second order system to describe gestures. Therefore, we need both positions and velocities for our residual error measurements (see below). Velocity can be obtained through the use of an estimator or by taking a difference of position measurements. Unfortunately, using differences adds noise to the data, which could make parameter identification difficult.
- 15 • Data order dependent versus independent (for on-line techniques). Certain on-line techniques will produce different parameter values based on the order the gesture data is presented. Because we define a gesture as a family of trajectories, with each trajectory in the same family equally valid, our method should be data order independent. In particular, different excursions through the same data set should
- 20 result in the same parameters at the end of the data acquisition phase.

- Linear versus Non-Linear. A model is a combination of linear and non-linear states and parameters. Although perfect (non human created) circular oscillatory motions can be described by a linear-in-parameters and linear-in-states model, a human created gesture may require a more complex model. Furthermore, our system can
5 recognize more complex oscillatory motions. Therefore, a method for identifying parameters in a richer non-linear model is needed, because non-linear models can represent a much broader range of motions.

We chose our gesture model and parameter determination scheme as follows. First, we decided to abandon off-line batch techniques in favor of on-line ones for reasons
10 already discussed above. The on-line method needs to be chosen carefully, because there are relatively few cases where it can be guaranteed that the estimated parameters will be equivalent to those resulting from off-line techniques applied to the entire data set as a whole.

Next, in an attempt to use only position data, we examined a Series-Parallel
15 Observer, which provides an *estimate* of the other unknown state for purely LIS and LIP systems. We abandoned this observer because it cannot adequately estimate parameters of non-perfect human gestures. Specifically, we could not extend the method to NLIS systems.

An on-line gradient descent method was examined, but for presently available
20 methods applicable to NLIP systems, there is no guarantee that the parameters will

converge towards their optimal values. Also, the parameters computed via this method are dependent on the order the data is presented.

A Linear Least Squares method (LLS) was examined next, which makes use of all the data independent of ordering. The resulting recursive LLS technique work for NLIP
5 models, and, therefore, allow us to examine more flexible and useful gesture models.
[See Cohen 96 for a detailed discussion of the reasons for the above decisions.]

The Recursive Linear Least Squares incrementally incorporates new data for determining the parameters which will best fit a set of data points to a given linear model. The recursive LLS method uses a tuning rule for updating the parameter vector θ without
10 inverting a matrix, creating a more computationally efficient LLS algorithm. A tuning rule is required, because each block of data will result in a different set of parameters, as illustrated in Figure 18. The separate graphs show that each pair of (x_i, y_i) data points results in a different best fitting θ line.

A method of incrementally updating the parameter θ is described below. For full
15 details, see [Kumar 86]. The concept is illustrated in Figure 19. After the first two data points determine the best fit line, each additional data point slightly adjusts the line to a new best fit. Each new data point will shift the line less and less due to the weighting auxiliary equation in the recursive LLS method. The formulation below describes how the weighting function operates.

The recursive (incremental) Linear Least Squares tuning method proceeds as follows. The tuning rule has the form:

$$\theta_{m+1} = g(x_m, \dot{x}_m, \theta_m)$$

- Suppose we have the output data \dot{x} and state data x up to time m , and from this data we
- 5 have already determined the best parameters θ for the set. From [Cohen 96] we know that at the next time step, with \dot{x}_{m+1} and x_{m+1} :

$$\theta_{m+1} = \left(\sum_{k=1}^{m+1} f_k^T f_k \right)^T \sum_{k=1}^{m+1} f_k^T \dot{x}_k$$

Define $R_{m+1} = \sum_{k=0}^{m+1} f_k^T f_k$. Then:

$$\begin{aligned} R_m &= \left(\sum_{k=0}^{m-1} f_k^T f_k \right)^{-1} + f_m^T f_m \\ &= R_{m-1} + f_m^T f_m \end{aligned}$$

- 10 which implies:

$$R_{m-1} = R_m - f_m^T f_m$$

Therefore:

$$\begin{aligned}
\theta_{m+1} &= R_{m+1}^{-1} \sum_{k=1}^{m+1} f_k^T \dot{x}_k \\
&= R_{m+1}^{-1} \left(\sum_{k=1}^m f_k^T \dot{x}_k + f_{m+1}^T \dot{x}_{m+1} \right) \\
&= R_{m+1}^{-1} \left(\sum_{k=1}^m f_k^T f_k \theta_m + f_{m+1}^T \dot{x}_{m+1} \right) \\
&= R_{m+1}^{-1} (R_m \theta_m + f_{m+1}^T \dot{x}_{m+1}) \\
&= R_{m+1}^{-1} ((R_{m+1} - f_{m+1}^T f_{m+1}) \theta_m + f_{m+1}^T \dot{x}_{m+1}) \\
&= R_{m+1}^{-1} (R_{m+1} \theta_m - f_{m+1}^T f_{m+1} \theta_m + f_{m+1}^T \dot{x}_{m+1}) \\
&= \theta_m - R_{m+1}^{-1} f_{m+1}^T (f_{m+1}^T \dot{x}_{m+1} - f_{m+1} \theta_m)
\end{aligned}$$

This is an update law for the R_{m+1} and θ_{m+1} terms. We still have to find the inverse of R_{m+1} at each time step. Fortunately, the matrix inversion lemma yields:

$$(R_m + f_m^T f_m)^{-1} = R_m^{-1} - R_m^{-1} f_m^T (f_m R_m^{-1} f_m^T + 1)^{-1} f_m R_m^{-1}$$

5 Therefore:

$$\begin{aligned}
R_{m+1}^{-1} &= (R_m + f_m^T f_m)^{-1} \\
&= R_m^{-1} - R_m^{-1} f_m^T (f_m R_m^{-1} f_m^T + 1)^{-1} f_m R_m^{-1}
\end{aligned}$$

The above equation is a recursive formula for R_{m+1}^{-1} that is not based on taking the inverse of a matrix. The initial value of R_0 is chosen as the identity matrix. See [Cohen 96] for a more formal discussion. If more importance is attached to recent data than to data

10 received in the remote past, then we can choose θ_m to minimize:

$$\sum_{k=0}^m \lambda^{m-k} (\dot{x}_k - f_k^T \theta_k)^2$$

where λ is termed the *forgetting factor* and is chosen with $0 < \lambda < 1$. This results in:

$$\begin{aligned}\theta_{m+1} &= \theta_m + R_{m+1}^{-1} f_{m+1}^T (\dot{x}_{m+1} - f_{m+1} \theta_m) \\ R_{m+1}^{-1} &= \frac{1}{\lambda} R_m^{-1} - \frac{1}{\lambda} R_m^{-1} f_m^T (f_m R_m^{-1} f_m^T + \lambda)^{-1} f_m R_m^{-1}\end{aligned}$$

The above recursive equation is the identifier in our gesture recognition system. This identifier allows us to represent gestures using a NLIP model, with the parameters
5 identified using an on-line computationally efficient data order independent technique.

We now determine the specific model used to represent oscillatory motion gestures.

Given that we modeled gestures using an LIP/NLIS representation, the following algorithm was used to determine the appropriate model (see Figure 20).

For the first step, we created phase-plane plots of the gestures to be modeled, as
10 illustrated in the last plot in Figure 15. A term in a differential equation model was composed of a parameter associated with combinations of multiplied state variables of various powers, that is, of the form $\theta_i x_1^j x_2^k$. An example model (of a one dimensional motion is):

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2\end{aligned}$$

15 Intuition was used to “guess” appropriate models that would best match the phase plane motions. Because we believed an acceleration model will sufficiently characterize the gestures in which we are interested, the \dot{x}_2 equation is the one modified with additional terms and parameters. For each model, the specific parameters for each gesture in the lexicon were computed using the LLS method (as discussed earlier).

The models were tested in simulation by measuring how well each tuned parameter model can predict the future states of its associated gesture (i.e. computing a total residual error). The model which best discriminates between gestures was chosen.

If none of the models can clearly discriminate between different gestures in a
5 lexicon, then new models are required. The heuristic we used was to add or delete specific terms, and determine if there was a significant change (good or bad) in the model's ability to discriminate gestures. Adding two specific terms to the above equation, that is, using the new model

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2 x_1^2 + \theta_4\end{aligned}$$

10 results in a model that is better able to discriminate between gestures.

The results of this process of modeling oscillating circles and lines are detailed in the remaining parts of this section. This process will also be used in the Phase II effort to determine an appropriate model to classify certain non-linear gestures.

A variety of linear-in-parameter models for good circle and line gesture
15 representations were tested. As before, each model represented only one dimension of motion, which was expanded to two or three for actually gesture recognition (i.e. an oscillating circle or line is formed when two or three of these decoupled models are present, one for each planar motion dimension). Again, x_1 is the position state, and x_2 is the velocity state. Five of these models are shown below. The determination of such

models illustrates how a new (and more comprehensive model) could be determined when required for more complex dynamic motions.

To use the models described here on a digital computer, a fourth-order Runge-Kutta integration method was used [Press 88]. Simulations showed that a sampling rate
5 of 10 Hz is sufficiently small to allow the use of this method.

The linear with offset component model is the most basic second order linear system. The offset component allows the model to represent gestures that are offset from the center of the image plane. It contains two parameters and is of the form:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2\end{aligned}$$

10

The Van der Pol equation is a slightly non-linear system, containing three parameters. The θ_2 and θ_3 parameters are attached to damping terms. This system is of the form:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2 x_1^2\end{aligned}$$

An offset component is added to the Van der Pol equation in this system. This
15 system has four parameters and is of the form:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2 x_1^2 + \theta_4\end{aligned}$$

A more non-linear system than the Van der Pol equations, the Higher Order Terms system contains additional spring-like components. This system has six parameters and is of the form:

$$\dot{x}_1 = x_2$$

$$\dot{x}_2 = \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 + \theta_4 x_2 + \theta_5 x_2 x_1^2 + \theta_6$$

The Velocity Damping Terms system has additional damping terms. It contains eight parameters and is of the form:

$$\dot{x}_1 = x_2$$

$$\dot{x}_2 = \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 + \theta_4 x_2 + \theta_5 x_2 x_1^2 + \theta_6 x_2^3 + \theta_7 x_1^2 x_2^3 + \theta_8$$

5 The use of simulations to determine the best gesture model for representing oscillating circles and lines is now detailed. We first details the residual measure calculation. Next the use of the residual measure to determine the best gesture model is described.

10 A predictor bin is composed of a model with parameters tuned to represent a specific gesture. The role of a bin is to determine a gesture's future position and velocity based on its current state. To measure the accuracy of the bin's prediction, we compared it to the next position and velocity of the gesture. The difference between the bin's prediction and the next gesture state is called the residual error. A bin predicting the future state of a gesture it represents will have a smaller residual error than a bin
15 predicting the future state of a gesture it does not represent.

The computation for the residual error is based on equation:

$$\dot{x}_k = F_k^T \theta$$

Recall that $f(x)$ is a two dimensional vector representing the gesture's position and velocity. Therefore \dot{x}_k is the gesture's velocity and acceleration at sample k . We

compute \dot{x}_k from the gesture's current and previous position and velocity. The parameter vector $\hat{\theta}$ is used to seed the predictor bin. Then:

$$\hat{x}_k = f_k^T \hat{\theta}$$

The residual error is then defined as the normalized difference between the actual
5 value of \dot{x}_k and the calculated value of \hat{x}_k :

$$res_err = \frac{\|\dot{x}_k - \hat{x}_k\|}{\|\dot{x}_k\|}$$

Figure 21 illustrates this concept. Consider the gesture at a given velocity and acceleration, sample k . At sample $k+1$, the predictions from each bin and the actual velocity and acceleration values are shown. The difference between a bin's predicted
10 values and the gesture's actual values (according to equation above) is the residual error for that particular bin.

The total residual error is the *res_err* summed for all data samples. The following section presents the residual calculation for each gesture with respect to each of the computed parameters.

15 We now detail how we determined which parameterization model for the predictor bin would best differentiate gestures. A data set of position and velocities of gestures is required to test each model. Using a vision system data was recorded for a slow, medium, and fast circular gesture. The data is the x and y position and velocity measurements from the image plane of the vision system, although for these simulations

only one of the dimensions is used. There is a small transition time when a human begins a gesture. This transient is usually less than a second long, but the residual error measurement has no meaning during this time. Therefore, gestures that last at least five seconds are used. The data recorded from a human gesture is termed "real gesture data."

5 The total residual error was calculated by subjecting each predictor bin to each gesture type. A measure of a model's usefulness is determined by examining the ratio of the lowest residual error to the next lowest residual error in each column. The worst "residual error ratio" is the smallest ratio from all the columns because it is easier to classify a gesture when the ratio is large.

10 The residual error results of the Linear with Offset Component are shown in Table 1. The residual errors for the slow and medium gestures, with respect to their associated bins, are an order of magnitude lower than the other errors in their columns. The residual error of the fast gesture, with respect to the fast gesture bin, is one-fourth the size of the closest residual error in its column (the medium gesture bin). Therefore, the Linear with
15 Offset Component system is a good candidate for a gesture model.

As seen in Table 2, the Van der Pol model is only a fair candidate for gesture discrimination. The residual error of the medium gesture with respect to its gesture bin is only two-fifths smaller than the residual error with respect to the slow gesture bin. Also, the residual errors in the slow gesture column are not an order of magnitude apart.

20 The Van der Pol with Offset Component model is better at discriminating gestures than the model without the offset term (see Table 3). The residual errors in the medium

gesture's column are now an order of magnitude apart. Although the residual errors in the fast gesture's column are not, the discrimination is still slightly better than in the Linear with Offset Component model.

Table 4 shows the residual errors associated with the Higher Order model. This model is an improvement over the Van der Pol with Offset Component model, as the residual errors in the fast gesture's column are now almost an order of magnitude apart.

Finally, Table 5 lists the residuals errors for the Velocity Damping model. This is the best model for discriminating between gestures, as the residual errors for each gesture with respect to their tuned bins are all at least an order of magnitude below the other residual errors in their columns.

A comparison of the worst "residual error ratio" of each model we considered is summarized in Figure 22, and suggests that the Velocity Damping model is the best choice for our application. However, the technique described here shows how more models could be derived and tested. For simple dynamic gesture applications, the Linear with Offset Component model would be used. For more complex gestures, a variation of the Velocity Damping model would be used.

Combining one dimensional motions to form higher dimensional gestures.

We have shown how predictors can be used to recognize one dimensional oscillatory motions. Recognition of higher dimensional motions is achieved by independently recognizing multiple, simultaneously created one dimensional motions.

For example, the combination of two oscillatory line motions performed in perpendicular axis can give rise to circular planar gestures, as shown in Figure 23.

Humans have the ability to create these planar motions. However, they can also make these motions in all three dimensions (for example, circles generated around different axis). To recognize these planar gestures, performed in three dimensional space, a vision system must be able to track a gesture's position through all three physical dimensions. A binocular vision system has this capability, as does a monocular system with an attached laser range finder. Any of these such vision systems can be used with our gesture recognition system to identify three dimensional gestures.

10

Development of a System to Recognize Static Gestures

Recognizing static hand gestures can be divided into localizing the hand from the rest of the image, describing the hand, and identifying that description. The module to recognize static hand gestures is to be both accurate and efficient. A time intensive process of evaluating hand gestures would prevent the system from updating and following motions which occur in real time. The system is intended to interact with people at a natural pace. Another important consideration is that the background may be cluttered with irrelevant objects. The algorithm should start at the hand and localize the hand from the surroundings.

20

Methodology

In order to meet these demands, the edges of the image are found with a Sobel operator. This is a very fast linear operation which finds approximations to the vertical and horizontal derivatives. In order to use only a single image, the greater of the horizontal and vertical component is kept as the value for each pixel. Besides being quick to calculate, an edge image avoids problems arising from attempting to define a region by locating consistent intensity values or even consistent changes in intensity. These values can vary dramatically in one hand and can be very hard to distinguish from the background as well.

In order to describe the hand, a box which tightly encloses the hand is first found. This allows a consistent description which is tolerant to changes in scale. To locate this box, we assume a point within the hand is given as a starting point. This is reasonable because the hand will be the principal moving object in the scene. Moving objects may be easily separated and the center of the largest moving area will be in the hand. From this starting point, a prospective box edge is drawn. If this box edge intersects an existing line, it must be expanded. Each side is tested in a spiral pattern of increasing radius from the initial center point. Once three sides have ceased expanding the last side is halted as well. Otherwise, the last side would often crawl up the length of the arm. The bounding box is shown in Figure 24.

Once the hand has been isolated with a bounding box, the hand is described (Figure 25). This description is meant to be scale invariant as the size of the hand can

vary in each camera image. At regular intervals along each edge the distance from the bounding edge to the hand's outline is measured. This provides a consistent description which may be rapidly calculated. A description is a vector of the measured distances, allowing a very concise representation.

- 5 The last task of the static gesture recognition is to identify the new description. A simple nearest neighbor metric is used to choose an identification. A file of recognized gestures is loaded in the initialization of the program. This file consists of a list of named gestures and their vector descriptions.

10 **Considerations**

- The primary obstacle in static gesture recognition is locating and separating the hand from the surroundings. Using sophisticated models of the hand or human body to identify with an image are computationally expensive. If orientation and scale are not very constrained, this cannot be done in real time. Our system makes descriptions
- 15 quickly and can compare them to predefined models quickly.

- The limitations of the current system are a result of being dependent on the fast edge finding techniques. If lighting is highly directional, parts of the hand may be placed in shadow. This can cause odd, irregular lines to be found and defeat the normal description. If the background immediately surrounding the hand is cluttered with
- 20 strongly contrasting areas, these unrelated lines may be grouped with the hand. This also causes unpredictable and unreliable descriptions. Such a background is very difficult to

separate without making assumptions about the hand color or the size of the hand. An upper and lower bound are placed on the size of the hand in the image, but these permit a wide range of distances to the camera and are needed to assure that enough of the hand exists on image to make a reasonable description.

- 5 As long as the hand is within the size bounds (more than a speck of three pixels and less than the entire field of view) and the immediate surroundings are fairly uniform, any hand gesture may be quickly and reliably recognized.

- Multiple camera views can be used to further refine the identification of static gestures. The best overall match from both views would be used to define and identify
10 the static gestures. Furthermore, the system works not just for “hand” gestures, but for any static type of gestures, including foot, limb, and full body gestures.

The Overall Gesture Recognition System

- In this section, based on the discussed functional and representational issues, we
15 detail the specific components of our dynamic gesture recognition system from an architectural and implementational viewpoint. Our system is composed of five modules: Figure 4 illustrates the signal flow of the gestural recognition and control system, from gesture creation, sensing, identification, and transformation into a system response.

- 20 **Gesture Creator**

In the Gesture Creator module, a human or device creates a spatial motion to be recognized by the sensor module. Our gesture recognition system was designed to recognize consistent yet non-perfect motion gestures and non-moving static gestures. Therefore, a human as well as a device can create the gestures which can be recognizable
5 by the system. Human gestures are more difficult to recognize due to the wide range of motions that humans recognize as the same gesture. We designed our gesture recognition system to recognize simple Lissajous gesture motions (repeating circles and lines), advanced motions such as “come here” and “go there”, and static hand symbols (such as “thumbs-up”).

10

Dynamic Gesture Lexicon

A gesture lexicon is a set of gestures used for communication or device control. We chose gestures for our lexicon based on the following:

- Humans should be able to make the gestures easily.
- 15 • Device gestures in the form of repeated motions should be modeled the same as human gestures.
- The gestures should be easily represented as a dynamic system.
- The lexicon should match useful gestures found in real world environments.

The dynamic gestures used in this system consist of three one-dimensional
20 oscillations, performed simultaneously in three dimensions (or two oscillations performed in two dimensions). A circle is such a motion, created by combining repeating motions in

two dimensions that have the same magnitude and frequency of oscillation, but with the individual motions 90 degrees out of phase. A “diagonal” line is another such motion. To illustrate this, we define three distinct circular gestures in terms of their frequency rates: slow, medium, and fast. Humans create gestures that we define as slow large
5 circles (slow), fast large circles (medium), and fast small circles (fast). More complex gestures can be generated and recognized, but these simple ones are used for illustrative purposes.

Main Three Gestures

10 Using the simpler Linear with Offset model (whose parameters are easier to understand than the more complex models), we represented a circle by two second order equations, one for each axis:

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= \theta_1 x_1 + \theta_2\end{aligned}$$

and

15

$$\begin{aligned}\dot{y}_1 &= y_2 \\ \dot{y}_2 &= \theta_1 y_1 + y_2\end{aligned}$$

Our gesture model has no “size” parameter. θ_1 is a frequency measure, and θ_2 is a drift component. The gestures were named “large”, “small”, “fast”, and “slow” due to the human motions used to determine the parameters (see Figure 26). A fast small circle is used to represent a fast oscillation because humans can not make fast oscillations using

large circles. Models with higher order terms would have parameters with different representations.

Expanded Lexicon - Geometric Constraints

5 A total of 24 gestures are possible from this example representation when the following are distinct gestures: clockwise and counter-clockwise circles, diagonal lines, one dimensional lines, and small and large circles and lines. Geometric constraints are required to expand the lexicon, because different gestures can result in the same parameters. Figure 5 shows motions that would cause an identifier to produce the same
10 frequency measure and drift components as it would produce when identifying a slow large circle. When x and y oscillating motions are 90 degrees out of phase, a clockwise circle is produced. Motions that are 270 degrees out of phase result in a counter clockwise circle. In phase motions produce a line with a positive slope. When the motions are 180 degrees out of phase, a line with a negative slope is produced. We can create additional
15 gestures from the fast small circle in the same manner.

 Given the various combinations of slow, fast, small, and large circles, the only one not used as a gesture is the slow small circle. Because, the slow small circle has the same oscillation frequency (medium) as the fast large circle, we need another geometric feature, the circle's size, to differentiate between these two gestures. As with the previous
20 gestures, additional gestures can be created from these two gestures by varying the phase relationships. Figure 7 shows a representation of the 24 gestures in this example lexicon.

Phase relationships are determined as follows. During the gesture, the x's and y's minimum and maximum image plane positions are computed. If the x and y motions are out of phase, as in a circle, then when x or y is minimum or maximum, the other axis's velocity is large. The clockwiseness of the motion is determined by looking at the sign of this velocity component. Similarly, if the x and y motion are in phase, then at these extremum points both velocities are small. A similar method is used when the gesture is performed in three dimensions.

Sensor Module

Unmodified Cohu Solid State CCD cameras are used as the sensor devices. No filters were used and the background was not modified. The Matrox Meteor capture card allows us to scale a captured image to any size without missing any frames. It will capture and transfer full-resolution, full-frame NTSC (640x480) or PAL (768x576) video input in real-time (30 Hz).

The color tracking system (CTS) uses the color of the hand and its motion to localize the hand in the scene. The hardware of the CTS system consists of a color camera, a frame grabber, and an IBM-PC compatible computer. The software consists of the image grabbing software and the tracking algorithm. Once the CTS is running, the graphical user interface displays the live image from the color camera on the computer monitor. The operator can then use the mouse to click on the hand in the image to select a

target for tracking. The system will then keep track of the moving target in the scene in real-time.

The color tracking system is developed on a BSD 4.0 UNIX operating system. The hardware involved consists of a color camera, an image capture board and an IBM
5 PC compatible. The software for the CTS is written in C and uses Motif for its graphical user interface (see Figure 27).

The present HTS system consists of a COHU 1322 color camera with a resolution of 494x768 pixels. The camera is connected to a Meteor image capturing board situated inside a Pentium-II 450MHz IBM-PC compatible computer. The Meteor board is capable
10 of capturing color video images at 30 frames per second. It is also able to capture these images at any resolution below the resolution of the camera.

The graphical user interface for the CTS displays a live color image from the camera on the computer screen. The user can then identify the target in the scene and click on it using the mouse. The CTS will then track the target in real-time. The flow
15 chart of the tracking algorithm is shown in Figure 28.

We capture the image using functions from the Meteor driver. To provide real-time operation, we setup the board to signal the program using a system interrupt (SIGUSR2). Every time a new frame is ready, the Meteor alerts the program with an interrupt on this signal. The image capture function responds to the interrupt by
20 transferring the current camera image to a buffer and processing it to find the target. The signal mechanism and its handling are what enable the system to operate in real-time.

The graphical user interface of CTS displays the live camera image on the screen. The user can start tracking by clicking the mouse on the target. This starts the tracking algorithm. The graphical user interface of the CTS is shown in Figure 29.

Once the user clicks on the target in the image, we compute the average color of a
5 small region around this point in the image. This will be the color of the target region being tracked in the scene until it is reinitialized. Once tracking begins, we compute the position of the target region in the image using two methods. The first method tracks the target when there is sufficient motion of the target in the image. The second method will take over when there is no motion of the target in the scene.

10 Before choosing the methods for finding the target in the scene, the system checks for motion in a region near the current or estimated target position using a motion detecting function. This function computes the difference between the current image and the previous image, which is stored in memory. If motion has occurred there will be sufficient change in the intensities in the region. This will indicate motion. The motion
15 detection function will trigger if a sufficient number of pixels change intensity by a certain threshold value.

If the motion detection function detects motion, the next step is to locate the target. This is done using the difference image and the target color. When an object moves between frames in a relatively stationary background, the color of the pixels
20 changes between frames near the target (unless the target and the background are of the same color). We compute the color change between frames for pixels near the target

location. The pixels whose color changes beyond a threshold make up the difference image. Note that the difference image will have areas, which are complementary. The pixels where the object used to be will complement those pixels where the object is at now. If we separate these pixels using the color of the target, we can compute the new
5 location of the target. The set of pixels in the difference image, which has the color of the target in the new image, will correspond to the leading edge of the target in the new image. If we assume that the target approximates an ellipse of known dimensions, we can compute the position of the center of the target (ellipse) from this difference image (see Figure 30).

10 The color of a pixel in a color image is determined by the values of the Red, Green and Blue bytes corresponding to the pixel in the image buffer. This color value will form a point in the three-dimensional RGB color space (see Figure 31). For our tracking system, when we compute the average color of the target, we assume that the target is fairly evenly colored and the illumination stays relatively the same. The average
15 color of the target is then the average RGB values of a sample set of pixels constituting the target. When the target moves and the illumination changes the color of the target is likely to change. The color matching function allows us to compute whether a pixel color matches the target color within limits. When the illumination on the target changes, the intensity of the color will change. This will appear as a movement along the RGB color
20 vector as shown in the figure below. In order to account for slight variations in the color, we further allow the point in color space to lie within a small-truncated cone as shown in

the figure. Two thresholds will decide the shape of the cone. One for the angle of the cone and one for the minimum length of the color vector. Thus, any pixel whose color lies within the truncated cone in color space will be considered as having the same color as the target.

5 When the motion detection function fails to detect significant motion in the scene, we use a static target matching function to compute its location. The function searches a small area about the current location of the target to find the best fit in the image for the target. The search will find the location of the target with the highest matching value. We assume that the object is approximately elliptical. The elliptical target is hypothesized at
10 each point in the search space and the matching metric is computed. This matching metric function uses a combination of edge and interior color matching algorithms to get a single matching number.

The image capture board is capable of providing us with a 480x640-pixel color image at 30 frames per second. Processing such a large image will slow down the
15 program. Fortunately, the nature of the tracking task is such that, only a fraction of the image is of interest. This region called the window of interest lies around the estimated position of the target in the new image. We can compute the location of the target in the new image from the location of the target in the previous image and its velocity. This simple method is able to keep track of the target even when the target moves rapidly. We
20 have found that the window of interest is typically $1/100^{\text{th}}$ the area of the original image. This speeds up the computation of the new target location considerably. A computer with

a higher processing speed could process the entire image without resorting to creating a region of interest.

Identification Module

5 The gesture recognition algorithms are located in the Identification Module. This module uses the position and velocity information provided by the sensor module to identify the gesture. The module is shown in Figure 32 and consists of three components - the Dynamic Gesture Prediction module, the Static Gesture Identification module, and the Overall Determination module (Which Gesture?). The output of the Overall
10 Determination module is sent to a screen display and to the SSM which produces an output based on the gesture command received.

The Dynamic Gesture Prediction Module

 The Dynamic Gesture Prediction module contains a bank of predictor bins (see
15 Figure 33). Each predictor bin contains a dynamic system model with parameters preset to a specific gesture. We assumed that the motions of human circular gestures are decoupled in x and y. Therefore, there are separate predictor bins for the x and y axes. In this example of three basic two dimensional gestures, a total of six predictor bins are required. The position and velocity information from the sensor module is fed directly
20 into each bin.

The idea for seeding each bin with different parameters was inspired by Narendra and Balakrishnan's work on improving the transient response of adaptive control system. In this work, they create a bank of indirect controllers which are tuned on line but whose identification models have different initial estimates of the plant parameters. When the
5 plant is identified, the bin that best matches that identification supplies a required control strategy for the system [Narendra 94].

Each bin's model, which has parameters that tune it to a specific gesture, is used to predict the future position and velocity of the motion. This prediction is made by feeding the current state of the motion into the gesture model. This prediction is compared to the
10 next position and velocity, and a residual error is computed. The bin, for each axis, with the least residual error is the best gesture match. If the best gesture match is not below a predefined threshold (which is a measure of how much variation from a specific gesture is allowed), then the result is ignored; no gesture is identified. Otherwise, geometric information is used to constrain the gesture further. A single gesture identification
15 number, which represents the combination of the best x bin, the best y bin, and the geometric information, is outputted to the transformation module. This number (or NULL if no gesture is identified) is outputted immediately upon the initiation of the gesture and is continually updated.

The parameters used to initially seed each predictor bin were calculated by
20 feeding the data of each axis from the three example basic gestures into the recursive linear least squares. The values for each bin are summarized in Table 6.

The Static Gesture Identification Module

The Static Gesture Identification module only searches for static gestures when the hand motion is very slow (i.e. the norm of the x and y velocities is below a threshold amount). When this happens, the module continually identifies a static gesture or outputs
5 that no gesture was found.

The static gestures may be easily expanded by writing new gesture descriptions to a configuration file. Each gesture is described by a name tag, width, height, x location, y location, base side, and three vectors (in this example, each consisting of 15 integers)
10 describing the profile of the hand. Because profiles may be significantly different due to varying tilts of the hand, multiple descriptions of fundamentally the same gesture may be desired. The initial or last line may also be less reliable due to missing the contours of the hand edge image.

The following are example parameter files (see Table 7). In each the name string
15 is followed by an arm side, width, height, x location and y location. The arm parameter is simply an integer corresponding to above, below, right, or left. The width and height are measured in pixels. The x and y location are 0 if the location is not important or +1 or -1 to restrict recognition of a gesture to one particular quadrant. The following three vectors are the extreme side (the end of the hand) then the top or left side followed by the bottom
20 or right side. The determination of which side is being represented is determined by the

arm side parameter. For example, if the base side is from below (as in the Halt gesture below) the first line is from above , then from the left , then from the right. Right and left refer to the overall image - not the facing of the imaged person.

Another method used for this part is to parameterize each part of the hand (palm, digits, and wrist) as a set of connected “blobs”, that is, three dimensional shapes which are connected together geometrically. As before, a configuration file would be used to define how these blobs are connected, with the vision system identifying the blobs which this module sticks together.

10 The Overall Determination Module

This “Which Gesture?” module takes input from both the Static and Dynamic Gesture modules. When the velocity is small, then a static gesture is observed. When the velocity is greater than a threshold value, then a dynamic gesture is observed. The gesture identified is continuously outputted, and can therefore change value over time (the output can even be that no gesture was identified). The gesture identified is sent to the transformation module.

Transformation Module

The transformation module take a gesture type as its input and determines what to do with it. In the case of this system, the gesture is converted to parameters which

represent the static or dynamic gesture, which is sent to the system which uses this information to produce a response.

System Response

5 The gesture command can be used for a wide variety of purposes. These include:

- Commands into a virtual reality simulator, to control and interact with the environment.
- Commands for a self service machine (SSM), such as a public information kiosk or Automated Teller Machines.
- 10 • Commands to control an actuated mechanism, such as a robot arm or mobile robot.
- Commands to control any device (such as a home appliance).

It is important to note that these devices can be controlled using static gestures, dynamic gestures, or a combination of the two. Thus, there is more information available to these system from the gesture input device, thereby allowing for a greater ability for
15 humans to command and control them.

The key features of our architecture are the prediction modules and the signal flow from gesture creation to system response. The other modules could be replaced with functionally equivalent systems without changing the structure of our architecture. For example, instead of a human, a robot could create the gesture. Alternatively, one could
20 create the gesture using a stylus, with a graphics tablet replacing the vision system in

sensor module S. The graphics tablet would output the x and y coordinates to the identification module I. Similarly, module R could be a robot, one as complex as a six degree of freedom robot arm or as simple as a stepper motor based camera platform. The former mechanism requires a more complex transformation scheme in module T, while
5 the latter system needs only a simple high level command generator.

As discussed earlier, the static and dynamic identification modules contains the majority of the required processing. Compared to most of the systems developed for gesture recognition (for example, see [Darrell 93] and [Murakami 91], for more details, see [Cohen 96]), this system requires relatively little processing time and memory to
10 identify one gesture feature. This makes it possible to create a system with the ability to identify multiple features in parallel. A sophisticated module could then examine the parallel gesture features and infer some higher level motion or command.

The Behavior Recognition System

15 Just as the gesture recognition module is built on a bank of predictor bins, the behavior recognition system is composed of a bank of gesture recognition modules. Each module focuses on a specific point of the body (such as a foot or knee). As that point moves through space, a "gesture" is generated and identified. The combination of gestures from those points are what we define as a motion behavior, which can be
20 categorized and identified. The system, illustrated in Figure 34, details the behavior

recognition system. The simplicity of the behavior recognition system is possible because of the strength and utility of the gesture recognition modules.

Overall System Flow

5 The signal flow proceeds as follows. A user is tagged at various body locations automatically by the vision system. The data is acquired at 30 Hz and sent to a parser which splits off the data from each specific body location to its own gesture recognition module (GRM). There is one GRM for each tagged feature. Each GRM outputs which gesture it recognized (if any) which is again sent to an identification
10 module. The identification module matches the gestures to their body location, defining a behavior. If this behavior matches one from a set of predefined behaviors, then this information is outputted.

The Parser

15 In the Parser module, the data, which is input as a stream of consecutive x,y,z,time coordinates from each tagged body location, is split up according to body location and sent to an appropriate GMI. This module needs to be changed whenever the input data is of a different format. Runtime variables define how many body parts are being tracked, and therefore the parser uses this information to determine the number of GMI bins and
20 how to split up the data properly.

Gesture Recognition Modules (GRMs)

The time and coordinate data from each body feature is used as inputs to an appropriate GRM. Each GRM module is exactly as described in section II a, except that these modules handle three dimensional points instead of just two.

5 Behavior Identification Module

The Behavior Identification Module accepts as inputs gesture types and body locations from the GRMs. Various combinations of gestures at specific body locations are designated as behaviors, and if a match is found, then the program outputs that match.

10 Behavior Recognition Experiments and Data

We performed experiments to test the behavior recognition system. First, a variety of behaviors were performed and served as the baseline for identification. Then these behaviors were repeated and the data sent through the system.

The behaviors centered on repeated leg and waist motion. The three main types of behaviors performed were:

Jumping Jacks: Start by standing straight, legs together, then move the feet out sideways and back together.

Walking in place: Move the feet in a slow walking motion.

Jogging in place: Move the feet up and down in an exaggerated quick motion.

Other behaviors used for experimentation, some of which included full three-dimensional motion:

Squats: Bend the knees to lower the body, then rise again.

Walking: Walk slowly towards the sensors.

Running: Run towards the sensors.

Skipping: Skip towards the sensors.

- 5 Hopscotch: Jump and land on one foot, then jump and land on both feet spread apart.

Rigorous experiments were performed using the first three behaviors. We used the Velocity Damping Terms gesture model to capture the richness of the presented motions. There was clear discrimination between the motions identified by the various
10 identification bins, and, as shown below, behavior recognition was possible.

For the Jumping Jacks behavior, the eight theta's for each axis for each sensor are:

- sensor 1:
x-axis: 26.8199, 57.4879, 0, 0, 0, 1751.95, 0, 0,
y-axis: 0, 0, -257.759, 0, 0, 15.9921, 58.561, 0,
15 z-axis: 0, 0, 24.4205, 57.9981, 0, 0, 0, -1026.36,
sensor 2:
x-axis: 17.8334, 58.4356, 0, 0, 0, 1691.1, 0, 0,
y-axis: 0, 0, -35.2513, 0, 0, 6.14527, 59.7759, 0,
z-axis: 0, 0, 28.6432, 57.4095, 0, 0, 0, 918.332,
20 sensor 3:
x-axis: 15.0551, 58.7612, 0, 0, 0, 1186.79, 0, 0,
y-axis: 0, 0, -476.275, 0, 0, 7.1385, 59.4896, 0,
z-axis: 0, 0, 6.74113, 59.9307, 0, 0, 0, -544.907
sensor 4:
25 x-axis: 10.8193, 58.9695, 0, 0, 0, -1210.42, 0, 0,
y-axis: 0, 0, 341.434, 0, 0, 9.92934, 59.1288, 0,
z-axis: 0, 0, 24.3505, 58.0358, 0, 0, 0, -142.869
sensor 5:
30 x-axis: 16.7313, 58.3415, 0, 0, 0, 4060.06, 0, 0,
y-axis: 0, 0, 819.198, 0, 0, 15.3747, 58.6259, 0,
z-axis: 0, 0, 27.1073, 57.975, 0, 0, 0, 612.659.

The residual results are shown in Table 8. As anticipated, the bold diagonal cells have lower residuals than all the other cells in their respective rows for each gesture in the behavior, demonstrating that our system can discriminate this behavior type. Some of the residuals in a bin are very close in value to others sensors. This is because they are, in fact, the same gesture. For example, sensor 4 and sensor 5 are placed on the knee and foot of the same leg, which would make the same oscillatory motion during a behavior, so of course the residual values would be similar.

For the Walking behavior, the eight theta's for each axis for each sensor are:

10 sensor 1:
x-axis: 24.304, 57.4632, 0, 0, 0, -351.126, 0, 0,
y-axis: 0, 0, -168.974, 0, 0, 23.2088, 57.6762, 0,
z-axis: 0, 0, 22.6047, 57.7623, 0, 0, 0, 1150.72,
sensor 2:
15 x-axis: 19.8496, 57.8372, 0, 0, 0, -1017.67, 0, 0,
y-axis: 0, 0, 31.8642, 0, 0, 26.8075, 57.5024, 0,
z-axis: 0, 0, 25.6103, 57.9468, 0, 0.000123365, 0, -358.633,
sensor 3:
20 x-axis: 10.2042, 59.182, 0, 0, 0, -617.508, 0, 0,
y-axis: 0, 0, 498.471, 0, 0, 30.6624, 56.95, 0,
z-axis: 0, 0, 22.534, 57.7824, 0, 0, 0, 598.156,
sensor 4:
25 x-axis: 24.6263, 57.6805, 0, 0, 0, -73.7837, 0, 0,
y-axis: -0.000250812, 0, 125.269, 0, 0, 19.2459, 58.205, 0,
z-axis: 0, 0, 18.6133, 58.2271, 0, 0, 0, -195.928,
sensor 5:
30 x-axis: 26.3399, 57.6672, 0, 0, 0, -522.552, 0, 0,
y-axis: -0.000136806, 0, 115.618, 0, 0, 18.8083, 58.326, 0,
z-axis: 0, 0, 19.0266, 58.1954, 0, 0, 0, 474.65,

The residual results are shown below in Table 9. Again, the bold diagonal cells have lower residuals than all the other cells in their respective rows for each gesture in the behavior, demonstrating that our system can discriminate this behavior type.

For the Running behavior, the eight theta's for each axis for each sensor are:

- 5 sensor 1:
x-axis: 22.9112, 57.8263, 0, 0, 0, -1763.57, 0, 0,
y-axis: 0, 0, -489.467, 0, 0, 11.7958, 58.779, 0,
z-axis: 0, 0, -3.51229, 61.0138, 0, 0, 0, 713.328,
sensor 2:
10 x-axis: -2.11517, 60.7157, 0, 0, 0, -40235.2, 0, 0,
y-axis: 0, 0, -4506, 0, 0, -79.0879, 70.5397, 0,
z-axis: 0, 0, -78.8084, 70.6087, 0, 0, 0, -375964,
sensor 3:
15 x-axis: 24.5412, 57.6338, 0, 0, 0, -2805.13, 0, 0,
y-axis: 0, 0, -211.096, 0, 0, 23.1457, 57.6718, 0,
z-axis: 0, 0, 20.9598, 58.3911, 0, 0, 0, 773.77
sensor 4:
20 x-axis: 20.1377, 58.218, 0, 0, 0, 4557.85, 0, 0,
y-axis: 0, 0, 607.713, 0, 0, 11.9292, 59.0339, 0,
z-axis: 0, 0, 16.2398, 58.6524, 0, 0, 0, -2667.72,
sensor 5:
25 x-axis: 29.6411, 56.9948, 0, 0, 0, 1093.19, 0, 0,
y-axis: 0, 0, 954.695, 0, 0, 14.4107, 58.6439, 0,
z-axis: 0, 0, 20.9606, 58.0327, 0, 0, 0, 3497.27.

- 25 The residual results are shown below in Table 10. As before, the bold diagonal cells have lower residuals than all the other cells in their respective rows for each gesture in the behavior, demonstrating that our system can discriminate this behavior type. As for quadruped locomotion, Kelso [Kelso 95] hints that there is dynamic mechanical
30 abstraction which can be studied to explain certain features of quadruped locomotion, specifically the limb frequencies of animals moving about the Serengeti plains. As he

states: "When plotted against limb length or mass, the stepping frequency, from Thompson's gazelle to the black rhinoceros, falls on three straight lines, one for each locomotory mode." Thus, it should be apparent to one of skill in the art that the invention is applicable to quadruped locomotion as well through appropriate extension.

5

Application Areas

Although one application of the invention is threat assessment, other uses are possible, including the generalized discrimination of suspicious or otherwise curious behaviors from normal activities. The system will not only perform surveillance activities, but make critical visual surveillance observations of locations for which there are not enough personnel to cover, or which are just simply too dangerous. Such surveillance tasks may include:

- Classifying: identify objects or classes of objects.
- Tracking: identify moving objects or people, and predict their future behavior.
- 15 • Patrolling: determine when something new enters or leaves the field of view, then track and classify it.
- Warning: determine threats in the area that need to be signaled back to an appropriate area. Such situations include, but are not limited to, vehicles speeding towards a facility, people loitering, climbing perimeter walls, carrying weapons, and transferring
20 objects.

In the commercial area, this system provides the foundation for a large number of gesture recognition and behavior recognition applications. Following is a list of example commercial applications to which the technology is applicable:

ATM Control. Given our understanding of the gesture recognition, the gesture
5 recognition system could be used for the control of Automatic Teller Machines (ATMs), as well as any other self-service machine or kiosk.

Crime Prevention. As discussed earlier, combinations of body gestures can be viewed as motion gaits, which can in turn be possibly identified as certain physical activities. Therefore, the same camera that is detecting gestures for the ATM (and any
10 other self-service machine or kiosk) control can be used to identify the physical activity in the area. Such activities could include possible robberies or assaults. Having such crime detection features in place would greatly increase the safety of using an ATM.

General User Interfaces. This gesture and behavior recognition system can be incorporated into general user interfaces, to augment other input devices such as the
15 mouse, keyboard, and voice commands. A system which could recognize behaviors would be useful in training people to use such devices.

As in the commercial arena, the gesture and behavior recognition system will have many uses in the military arena. The following is a list of military applications which we can expect to pursue in future projects.

Enhanced User Interfaces. As in the commercial application, the gesture and behavior recognition system can be incorporated into military computer systems that would benefit from an additional gestural user interface.

Surveillance. Since the dynamic gesture recognition can be used to detect
5 physical activity at an ATM, it can also be used to detect activities in a combat environment. These activities not only include running, walking, and firing a weapon, but would also include the detection of group behaviors and formations.

We claim: